

Interpretation of Handwritten Addresses in US Mailstream

Sargur N. Srihari, Venu Govindaraju and Ajay Shekhawat
Center of Excellence for Document Analysis and Recognition (CEDAR)
State University of New York at Buffalo
UB Commons, Suite 202
Amherst, NY 14228-2567, USA

Abstract

Determining the delivery location for mail pieces based on handwritten addresses is a problem that trained humans can normally solve. As a problem in machine reading and interpretation, it presents many challenges. A method for determining the delivery point codes (DPCs) for handwritten addresses by computer is described. Solution to the task requires locating and recognizing address components (e.g., ZIP Code, street number, P.O. box number) and using multiple information sources to assign the DPC to an address. Our methodology uses diverse pattern recognition techniques, image processing algorithms (thresholding, underline removal, separation of lines, location and recognition of address components), and access to United States Postal Service (USPS) databases to determine the DPC.

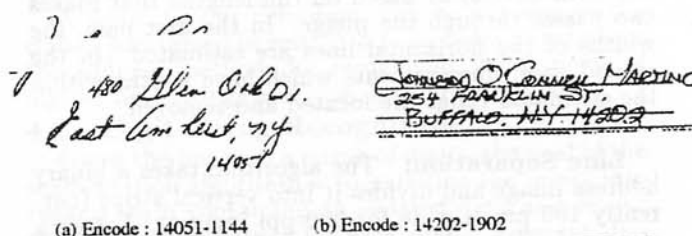
1 Introduction

The United States Postal Service sorts mail by utilizing an encoding of each mail piece by a string of digits, the first five of which are referred to as the ZIP Code, and together with the next four (add-on) as the ZIP+4. The first three digits encode the state, city and postal sectional center facility (SCF), the next two the post office, and the final digits encode destination information such as block face, post office box number, etc. The destination encoding can be a five-digit ZIP Code or nine-digit ZIP Code.

Examples of address blocks taken from live mail (Figure 1) illustrate some of the challenges posed by unconstrained handwriting and wide variations in address structure. Figure 1a is a very poorly handwritten address block in which the ZIP Code is located in a different line from the city name and the state name. The normal syntax would have the city name, state name, and ZIP all on the last line. Descenders from higher lines intrude into lines below. It is often the case, as in Figure 1a, that a straight line cannot be found which separates the adjacent lines without cutting through component(s). Figure 1b is an example of an address block with guide lines which must be treated as noise by the address interpreter.

A complete system for unconstrained address interpretation needs several components. Inputs to the system are a binarized image of the destination address block and USPS databases to "look-up" the add-on after the ZIP code and street address/ PO Box address

Figure 1: Examples (a) abnormal address syntax, (b) interfering guidelines



from the image are matched against a record in the database. Output is the DPC encoding (ZIP+add-on) for the mail piece. Figure 1 shows the DPC output obtained by the system on processing the addresses shown. At the pre-processing level, the written matter from the handwritten address block image must be distinguished from the background, and irrelevant streaks such as underlines must be removed. The text line locations in the address must be determined and the text lines separated into words. The ZIP Code, street number/P.O. Box number etc. must be read and the DPC determined.

There are three methods of determining the DPC as described below [7]:

1. The ZIP code happens to be a unique ZIP as identified by the City-State File¹. In such cases, the 5-digit ZIP Code itself is the DPC.
2. The address has a P.O. Box number. In such cases, the 5-digit ZIP Code and the P.O. Box number are sufficient to determine the 4-digit add-on with a ZIP+4 File² lookup. The 9-digit num-

¹City-State File: ZIP codes are associated with their city and county names. This directory has over 100,000 records. The records define, among other items, the type of ZIP code. For instance, "type 90" identifies unique ZIP codes where the 5-digit ZIP is itself the final destination point code, "type 14" identifies ZIP codes which have only PO Box addresses, etc.

²ZIP+4 File: It has 28 million records, where each record contains, among other items, ZIP code, add-on, record type (street, high-rise, PO Box, firm, etc.), carrier route ID, and street name information.

ber becomes the DPC.

3. The address has a street line (*street number + street name*). In such cases, the ZIP Code and the street number confine the possible street names to a small set with a DSF³ lookup. The word recognition algorithm is invoked to find which of the street names in the set matches the street name image in the address. The selected street name has a 4 digit add-on which can be obtained from the ZIP+4 File.

2 Preprocessing

Underline Removal: Approximately 5% of the time, the words in a handwritten address are written on machine printed guide lines that are provided to assist the writer. Often the handwritten text intersects these guide lines. An underline removal method has been developed based on run-lengths that makes two passes through the image. In the first pass, the widths of the horizontal lines are estimated. In the second pass, line segments which have widths within the estimated range are located and removed.

Line Separation: The algorithm takes a binary address image and divides it into vertical strips (currently 100 pixels wide for 300 ppi images). For each strip, a horizontal projection profile is computed. A group of adjacent profile positions which have non-zero values form a block. Blocks which are horizontally adjacent and vertically overlapping are connected to form address lines. Heuristics are used to split and/or join blocks. Connected components that are located entirely in blocks from a single line are assigned to that line. Connected components that span more than one text line (Figure 1a) are split into two or more components and placed in separate text lines [6].

The shape and size of the connected components are used to refine the line segmentation. For instance, a short horizontal bar may be identified as having the shape of a dash or a disconnected 5-hat⁴. The identity of the horizontal bar is verified by confirming the presence of another component to its left which must be a digit. The constraints for shape and position of components are not disjoint. For instance, if a component is in a likely 'comma location', shape constraints are relaxed to identify the components as a comma regardless of its shape.

The current module is able to perform correct line-separation of addresses at a rate of about 93%. The module was tested on 2000 address images. An error is counted whenever the line assignment of components affects the interpretation of the address.

³**Delivery Sequence File (DSF):** It has a 100 million address information records where each record corresponds to a mail stop, whether it is a house, an apartment, or a suite in a building. Once the address on the mailpiece is matched to a DSF address record, the mail can be routed to the correct carrier and sorted into the carrier's walk sequence.

⁴Handwritten 5's are often written in two strokes. When the horizontal stroke is not connected to the other stroke, it is considered to be a disconnected 5-hat.

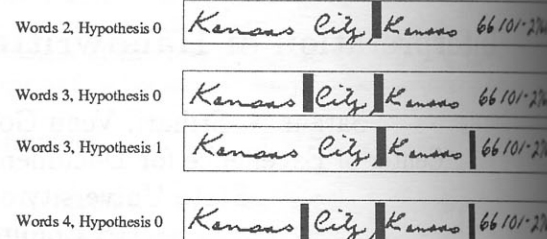


Figure 2: Word break hypotheses

Word Separation: Words are formed by ordering components of a text line from left to right and selecting word breaks between components. Word breaks are hypothesized based on distance between components, a shift in the mean vertical position of components, or the locations of punctuation (Figure 2).

Each text line can be divided into two to four words, and multiple hypotheses can be created for each number of words, i.e., a text line can be divided into three words in several different ways, depending on the location of the hypothesized word break in the line (Figure 2 shows two such hypothesis to obtain three words from the line). Each hypothesis is associated with a confidence ($Conf(Hypo_0) \geq Conf(Hypo_1) \geq Conf(Hypo_2) \dots$).

Finding a metric to compute the distance between adjacent components that conforms to human evaluation is non-trivial. Following measures have been used to estimate the distance between components: i) the horizontal distance between the bounding boxes of adjacent components, ii) the minimum horizontal run-length between adjacent components, and iii) the minimum Euclidean distance between adjacent components.

3 Parsing

The goal of parsing is to locate the ZIP code, street number + street name or "PO Box" + PO Box number in the address block. Input to the system is a list of connected components. Associated with each component is its line number assignment, coordinates, size, and the response of a digit recognizer. At this stage, text lines have been segmented into clusters of connected components to form a word.

Parsing proceeds at two levels. First, words are classified into categories based on the shape, size, and digit recognition response alone. Next, the classification is refined using syntax models for different lines of the address.

3.1 Shape Based Classification

Words are classified into one of the following categories: *text*, *city*, *state-abbreviation*, *digit-string* (along with an estimate of the number of digits in the string), *digits-dash*, *9-digit ZIP+4 Codes*, *barcodes*, *P.O. Box* and *noise*. Each word is identified as a dash,

Module	Accept rate	Correct rate
Street Number	91%	83%
PO Box number	70%	82%

Table 1: Street Number & PO Box number location and recognition performance

a comma, a disconnected 5-hat or none of these. In addition, each word is identified as a digit or a non-digit. There may be multiple identifications, e.g., a comma and the digit 1 may have identical shapes. The multiple identifications are reduced after the syntax information is considered [2].

3.2 Syntax-based Classification

Based on statistical analysis of information on mail pieces [5], the likely syntax constructions to be found in handwritten addresses is determined. Matches between the syntax, word groupings and word classifications are determined and the most consistent matches are examined in order to find likely ZIP code, street line, and/or PO Box candidates.

ZIP Code Recognition: Although the most likely syntax of the address has the bottom line containing the city, state, and ZIP Code, many other address syntaxes are possible. For instance, the bottom line may contain the ZIP Code and the second last line may contain the city and state name (Figure 1a). Words that match the syntax of a ZIP Code are selected as ZIP Code candidates. The system processes up to 5 candidates to determine the ZIP Code of the address (on average 1.6 candidates per mailpiece). The correct ZIP code is found to be present among the candidates in 83% of the cases. The correct candidate is chosen by using collaborative information from the other address components (e.g., confirm the ZIP code for compatibility with city name).

Street Number and P.O. Box Recognition:

Based on statistical analysis of address images, we know that the street address line and PO Box address line must be located beneath the top address text line and above the text line containing the city name. The street address line is located by searching for a line that contains digits followed by text. PO Box address line is located by searching for key words "PO", "Box", and periods. A digit recognizer response of "90" usually indicates the presence of "PO".

In a test on 2000 address images, there were 1465 images with street names and 446 images with PO Boxes⁵. The performance of the street number and PO Box modules is given in Table 1.

3.3 Reliability and confidence measures

If all digits are identified with sufficient confidence, the ZIP Code/ Street number/ PO Box number can-

⁵Some address images have neither street address nor PO Box address.

Data set	Length	Correct	rate
5863	1	5600	96%
620	2	486	78%
77	3	49	64%
8	4	2	25%

Table 2: Digit string recognition

didate is assigned a confidence based on recognition confidence of the individual digits and the amount of segmentation required. Presence of periods and digit recognition response of "90" boosts the confidence of a PO Box candidate. Presence of a string of nine digits with a dash boosts the confidence of a ZIP candidate. Presence of the # symbol on the right side of the line boosts the confidence of a street number candidate.

4 Handwritten Word Recognition

A word, by our definition, can be any one of the classes mentioned in section 3.1. In order to interpret the address, methods for recognizing words made of digits (digit-string) and alphabets are required.

4.1 Digit-string Recognition

Given the image of a string of digits, the goal of the segmentation algorithm is to partition the image into regions, each containing an isolated digit and return the recognition output for each digit. Adjacent digits can be touching and some of the digits might be broken into more than one component (e.g., 5-hats). Therefore the number of digits in a digit string is not simply a count of the number of components in the field.

The digit string recognizer can be invoked in three different modes: (i) estimate the number of digits (ii) force either 5 or 9 digits in the digit string and (iii) force a given digit string length [4]. The mode of estimating the number of digits is appropriate when acting on the street number and PO Box number fields, while the mode of specifying the length of the digit string (as 5 or 9) is appropriate when acting on the ZIP code field.

The segmenter has a correct segmentation rate of 92.9% when the input is specified as a 5 or 9 digit ZIP Code, and is 87.5% correct when the number of digits has to be estimated. Table 2 shows the performance of the digit-string recognizer falling as the number of touching digits in a single connected component increases [8].

4.2 Word Recognition

Lexicon	10	100	1000
HGR	95.8%	89.8%	78.8%
HMM	93.2%	84.0%	68.0%
Combination	97.4%	92.8%	84.4%

Table 3: Word Recognition performance on varying lexicon sizes

If the address has a street line (Street number + street name) and the street number is recognized with

Thresholds	% Encode	% Error
T ₁	43.6	5.4
T ₂	41.8	4.8
T ₃	40.8	4.2
T ₄	37.5	2.6
T ₅	31.2	1.6
T ₅	30.0	1.4

Table 4: Performance at several operating points

a high confidence, a dictionary of possible street names is generated from the DSF. The list of possible street names confined by the street number and a ZIP code is on average less than thirty. The dictionary and the street name image are passed on to the handwritten word recognition algorithms. Two approaches to word recognition have been explored.

One method is based on segmentation and recognition of individual cursive characters in the word [3]. The algorithm follows a *Hypothesis Generate and Reduce (HGR)* paradigm where the lexicon is continually reduced based on the word identity hypothesis at different stages of the algorithm. The segmentation process is actually a search that generates a series of segmentations of the word. Potential characters that result from segmentation are passed to a character recognizer. This results in a series of strings which represent the characters detected from each segmentation sequence of the word. The strings are passed to a dictionary matching routine which generates a ranking of the dictionary.

The second approach first divides the word into segments and then recognizes each segment in the string [1]. The recognition method is based on the *Hidden Markov Model (HMM)*. The observed segments are referred to as symbols. The characters giving rise to the symbols are the hidden states. The Viterbi algorithm and a dictionary are used to rank the dictionary.

The two algorithms (HGR & HMM) are combined to improve overall performance. Both methods were tested on 3000 images of handwritten postal words (state names, street names, etc.). Table 3 shows the performance of the individual algorithms and their combination.

5 Performance Evaluation

The system performance was tested on an unseen set of 1600 address images randomly picked from the mail stream. The performance at various error rates is shown in Figure 4. Encode rate corresponds to the % of the total number of images for which DPC was determined. An analysis of the errors indicates that about 20% of the errors are due to incorrect ZIP code recognition; about 16% of the errors are due to incorrect PO Box recognition; and the remaining are due to incorrect street number/street name recognition.

6 Summary

Handwritten Address interpretation is a challenging problem. It involves several document analysis techniques: i) binarization, ii) noise (postal marks,

underlines, etc.) removal, iii) skew correction, iv) address lines separation, v) identification of word categories (e.g., ZIP field, street number field), vi) recognition of words and digit strings, and vii) determining the destination point code (DPC) of the address. We have described an end-to-end system which integrates all these document analysis techniques to bear upon the task of address interpretation. The current performance is at 44% encode rate with $\leq 6\%$ error rate. The target is 50% encode rate with $\leq 2\%$ error rate.

References

1. M. Chen, A. Kundu and J. Zhou. Off-line handwritten word recognition using a single contextual Hidden Markov Model. *Proc. of the Computer Vision and Pattern Recognition Conf.*, Champaign, IL, pp. 669-672, 1992.
2. E. Cohen, J.J. Hull and S.N. Srihari. Understanding handwritten text in a structured environment: determining ZIP codes from addresses. *Int. Journal of Pattern Recognition and Artificial Intelligence*, vol. 5, no. 1&2, pp. 221-264, 1991.
3. J.T. Favata and S.N. Srihari. Off-line recognition of handwritten cursive words. *Proc. of the SPIE symposium on electronic imaging science and technology*, San Jose, CA, 1992.
4. R. Fenrich. Segmentation of automatically located handwritten words. *Proc. Int. workshop in handwriting recognition*, Bonas, FRANCE, pp. 33-44, 1991.
5. J.J. Hull, D.S. Lee and S.N. Srihari. Characteristics of handwritten mail addresses: A statistical study for developing an automatic ZIP Code recognition system. *TR 88-06, Dept. of Computer Science, SUNY Buffalo*, 1988.
6. S.N. Srihari, E. Cohen, J.J. Hull and L. Kuan. A system to locate and recognize ZIP Codes in handwritten addresses. *Int. Journal of Research and Engineering Postal Applications 1*, pp. 37-45, 1989.
7. S.N. Srihari, E. Cohen, V. Govindaraju, and A. Shekhawat. Determining the Delivery Point Codes on Handwritten Addresses. *Fifth Advanced Technology USPS Conference*, pp. 321-336, 1992.
8. G. Srikanthan, A. Shekhawat, and J.T. Favata. Gradient, Structural, and Concavity Features for Handwritten Digit Recognition. *In preparation*, 1993.